

Machine Learning

(main)

Supervised ML

Regression

- linear
- Ridge & Lasso
- polynomial Reg.
- Multi linear
- Non-linear
- OLS
- Time series forecasting

Classification

- logistic
- KNN
- DT
- Naive Bayes
- RF
- XGB
- GBM
- SVM
- ANN
- calboost, adaboost

Clustering

- kmeans
- DBScan
- hierarchical

Unsupervised ML

apriori

- apriori
- Recommender System
 - Content based
 - Memory based
 - User based
 - item based

Dimensionality Reduction

- PCA
- tSNE
- ICA
- SVD
- Auto

Encoding

Supervised ML →

We will start with Regression → linear

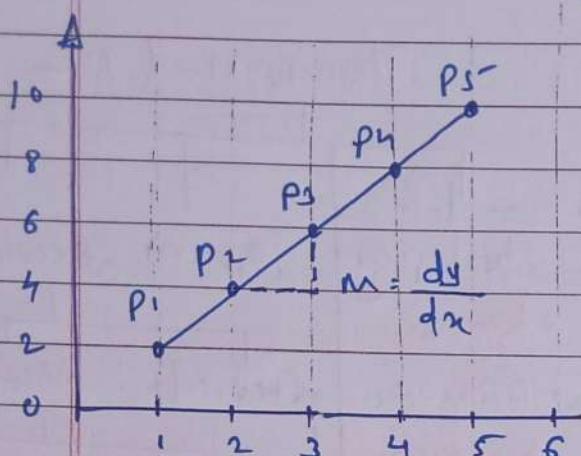
linear Regression

- Simple linear Regression (SLR)
- Multiple linear Regression (MLR)
- polynomial Regression

- Simple linear Regression (SLR) :- when for continuous output we have only one feature then it is called SLR

for Ex . CGPA | package
 6.0 2 LPA
 7.2 3 LPA
 8.9 6 LPA.

Basics of linear Regression.



Equation of straight line

$$y = mx + c$$

find value of m and c

$$P_2(2, 4) \rightarrow x_1, y_1$$

$$P_3(3, 6) \rightarrow x_2, y_2$$

$$m = \text{slope} = \frac{y_2 - y_1}{x_2 - x_1} = \frac{6 - 4}{3 - 2}$$

$$m = \frac{y_1}{x_1} = \frac{4}{2} = 2$$

slope means with every change in value of x how much differs in y

Intercept = c

$$P_4(x_3, y_3) \rightarrow (4, 8)$$

$$y = mx + c$$

$$8 = 2(4) + c$$

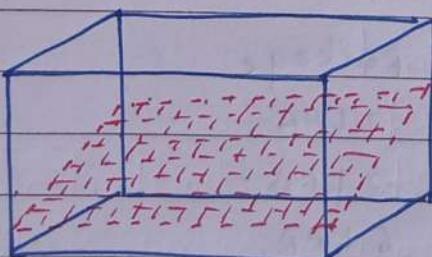
$$8 - 8 = c = 0$$

Inference : The above line eqn is function that relates x and y .

for given value of x we can find corresponding value of y .

what if we have two or more than two independent variables ?

→ then it is called multiple linear regression



(3D)

simple linear regression : $y = mx + c$

multiple linear regression.

$$y = a_1x_1 + a_2x_2 + a_3x_3 + \dots + a_nx_n + b.$$

Advantages : 1) simple to implement

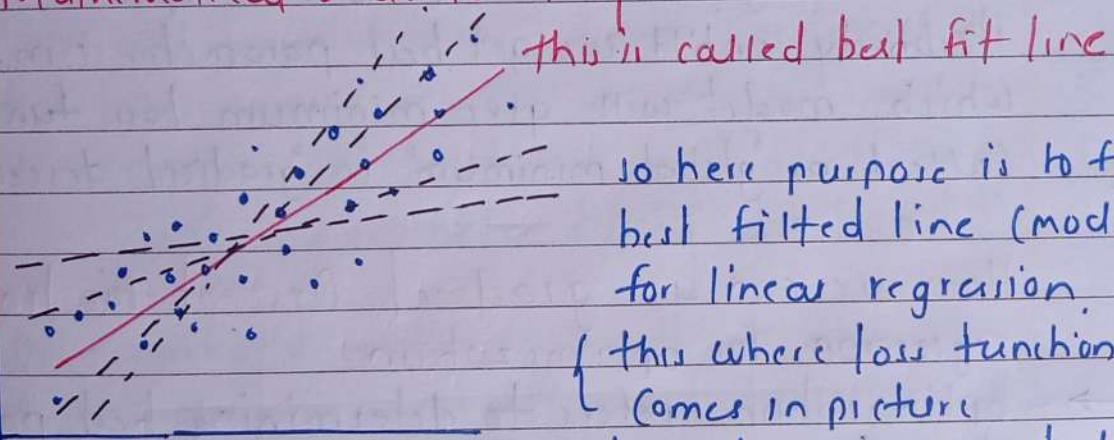
2) perform well on the data with linear relationship.

Disadvantages : 1) Not suitable for data having non-linear relationship.

2) underfitting issue

3) sensitive to outliers

Mathematical Understanding :-



so here purpose is to find best fitted line (model) for linear regression.

{ this where loss function comes in picture }

- loss function measures how far an estimated value from its true value.
- it is helpful to determine which model performs better and which parameters are better. (m, c)

$$\text{loss} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Let try to understand with the Ex. $m = 3, c = 2$

$$\hat{y} = 3x + 2$$

x	4	4
2	10	8
3	14	11
4	18	14
5	22	17
6	26	20

Now find its loss function.

$$\text{Loss} = \frac{[(10-8)^2 + (14-11)^2 + (18-14)^2 + (22-17)^2 + (26-20)^2]}{5}$$
$$= \frac{4+9+16+25+36}{5} = \frac{90}{5} = 18$$

{ reason to take square: so that positive and negative value may not cancel each other }

low loss value \rightarrow High Accuracy

high loss value \rightarrow low Accuracy

We can improve the model by some optimization technique called as "gradient descent" where repeat the process iteratively until we get best parameter (m, c) for which model will give minimum loss function. Called as "global minimum" in "gradient descent"

How we can use gradient Descent for linear regression for optimization.

→ Optimization refers to determining best parameters for model such that loss function of the model decreases as result of which model can predict more accurately.

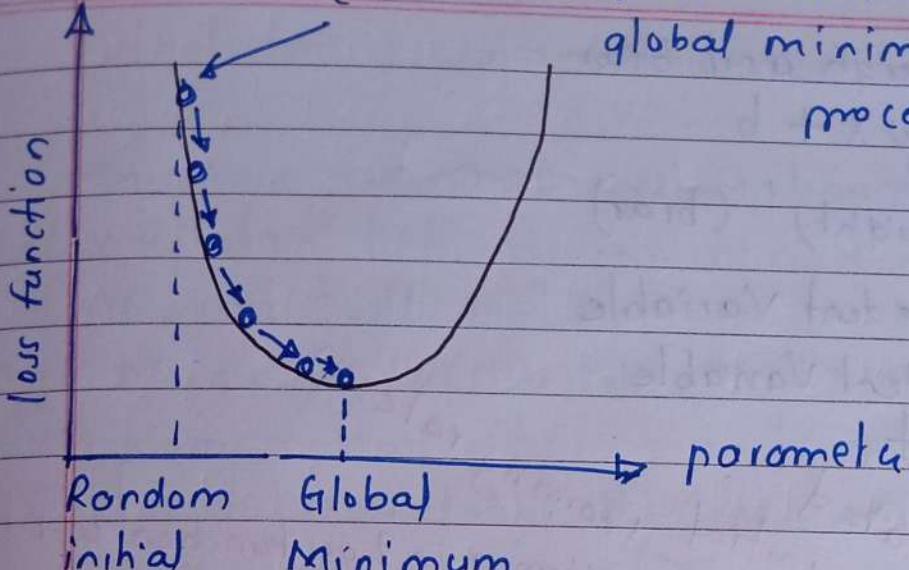
$$A \quad y = m_2x + c_2 \rightarrow \quad y = m_3x + c_3$$

$$y = m_1x + c_1$$

{ m and c are the parameters }

here we can find model 3 is best fit since the loss function is least and thus this model is optimum this where we use gradient descent to optimize model.

{this will step down until reach down to global minimum by iterative process}



Value
(m_1, c_1)

(m_3, c_3)

In machine learning

updated. $\rightarrow m = m_1 - LDm$ | $w = w - Ldw$
 $c = c_1 - LDc$ | $b = b - Ldb$

m - slope initial w = weight
 c - Intercept b = bias.

L - learning rate : it is magnitude of change that you want in parameter during iteration.

Dm : partial derivative of loss function w.r.t m

Dc : partial derivative of loss function w.r.t c .

$$Dm = \frac{\partial (\text{loss function})}{\partial m} = \frac{\partial}{\partial m} \left(\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \right)$$

$$= -2 \sum_{i=0}^n x_i (y_i - \hat{y}_i)$$

$$\text{Ily } Dc = \frac{\partial (\text{loss function})}{\partial c} = \frac{\partial}{\partial c} \left(\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \right)$$

$$= -2 \sum_{i=0}^n (y_i - \hat{y}_i)$$

In terms weight and bias

$$y = w x + b$$

↑
(weight) (bias)

x - independent Variable

y - dependent Variable

w - weight

b - bias

$$w = w - \alpha d_w$$

updated *initial learning rate (α)* *change in Loss function w.r.t w*

$$b = b - \alpha d_b$$

change in loss function w.r.t b

Learning rate :- it is tuning parameter in an optimization algorithm that determines step size at each iteration while moving toward minimum of loss function.

$$d_w = \frac{-2}{n} \sum_{i=0}^n x_i (y_i - \hat{y}_i)$$

$$d_b = \frac{-2}{n} \sum_{i=0}^n (y_i - \hat{y}_i).$$

for Multiple linear Regression. for one target we have two or more than two independent variables.

Prediction Equation for MLR

$$\hat{y} = q_1 x_1 + q_2 x_2 + \dots + q_n x_n + b.$$

and,

Regression Equation

$$y = q_1 x_1 + q_2 x_2 + q_3 x_3 + \dots + q_n x_n + b + \epsilon$$

Actual
Value

Predicted value

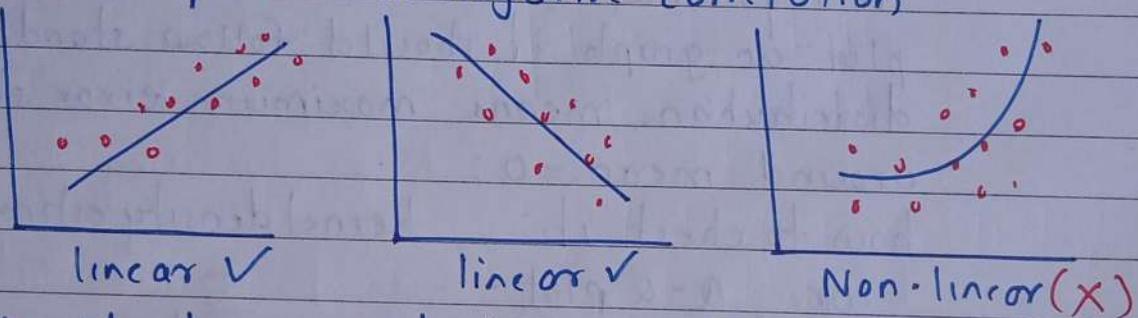
Error

What are assumptions of linear Regression?

there are five main assumption in linear regression.

1. Linear Relation between input and output.
2. No multicollinearity.
3. Normality of Residual.
4. Homoscedasticity.
5. No autocorrelation of error

1. Assumption 1:- there should be linear relation between individual feature and target (output) it could be positive or negative correlation



applicable to multiple linear problem as well. Notok
how to check this : scatter plot : (feature 1 Vr target)
(feature 2 Vr target).

2. Assumption 2:- Multicollinearity .. it mean there all the feature should be independent or should not have any correlation among themselves.

why , what the problem?

In multi linear Regression Model for 3D we draw a hyperplane.

$$y = q_1x_1 + q_2x_2 + q_3x_3 + b.$$

where q_i represent what will be the change in y with respect to x_1 assuming x_2 and x_3 constant. but if it violates this assumption then model will not perform good. (Ex of two physcs scientist)

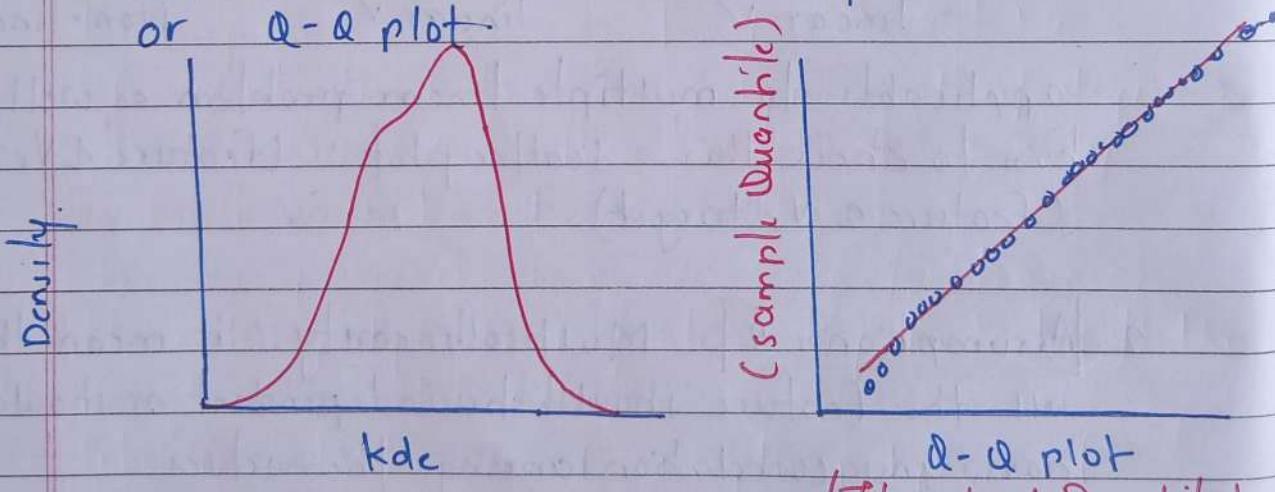
how to check multicollinearity :-

- 1) VIF (Variance Inflation factor) :-
if it is around 1 then features don't have the issue and if it is 5 or more than that then that particular feature has multicollinearity issue and need to remove it.
- 2). another method is to find out correlation between all the features (Heatmap)

3 Assumption 3 - Normality of Residual.

it says that when error (actual-predicted) plot or graph it should follow standard normal distribution means maximum error should around mean = 0.

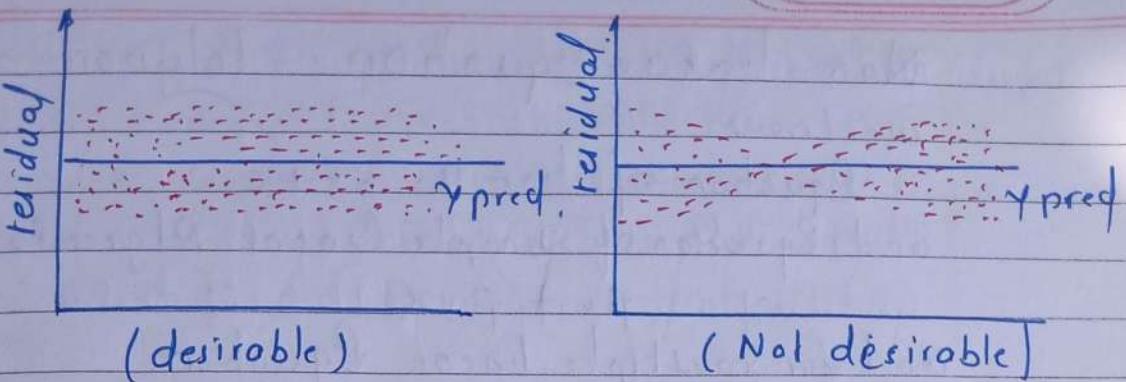
how to check it : kernel density estimation (kde)
or Q-Q plot



4 Assumption 4: Homoscedasticity.

some scattered (spred).

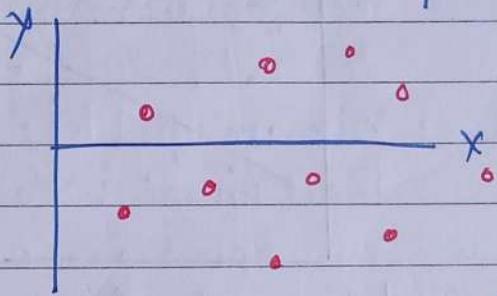
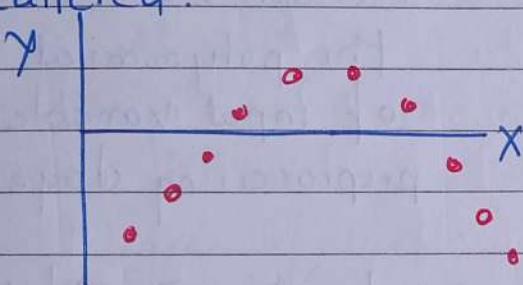
spread of residual should be equally or uniformly scattered. if it is not it called as heteroscedasticity which is not desirable.



how to check:- scatter plot (y_{pred}, residual).

Assumption 5: No autocorrelation of error.

If we plot the error there should not be any specific pattern instead it should randomly scattered.



positive Autocorrelation
(Not desirable)

Negative Autocorrelation
(desirable)

how to check : pli. plot (residual)

Summary chart

	Assumption	Severity	Prediction	Inference
1) Linear Relation	High	✓	✓	
2) Multicollinearity	Medium	✗		✓
3) Normality	Low	✗		✓
4) Homoscedasticity	High	✓		✓
5) Autocorrelation of Errors	-	-	-	-

Non-linear Equation - Polynomial Regression

We know

Equation of line $y = mx + c$

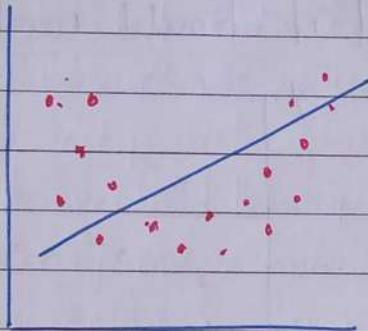
and Equation of simple Linear Regression

$$y = \beta_0 + \beta_1 x$$

and for multiple linear Equation.

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \dots + \beta_n x_n$$

This is applicable only when data is linear but what if data is not linear?



In such scenario we extract the polynomial features of input variable in preprocessing stage

Let say for ex $x | y \rightarrow 5 | 10$

- For x we want to make polynomial of degree 2 then we will convert $x \rightarrow x^0, x^1, x^2, y$
so it will be $1, 5, 25$
- this way we create a new data for training the extra polynomial feature try to extract this non-linear relationship.
- its formula become for simple polynomial regression.

$$y = \beta_0 + \beta_1 x + \beta_2 x^2$$

For degree 3

$$y = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3$$

* Now how we will know the perfect value for the degree — since this is hyperparameter.

if we keep it low then may be it cause underfitting mean, it may be not able to learn the all attributs. and if we select very high then there is chance of overfitting or overlearned that's why our job is to find out optimum value

- In case if we have two features x_1, x_2, Y then for degree 2 our simple polynomial Equation would become

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_1^2 + \beta_3 x_2 + \beta_4 x_2^2$$

Interview Ques: why polynomial Eqⁿ is called as Linear Regression

Ans: when we talk about linear Regression we talk about relation between y and coefficients of features and degree of coefficient is still one and thus relation between y and coefficient is still linear

Ordinary least square :- (OLS Algorithm)

- it is method for estimating the parameters of linear regression model.
- its aim to find the values of the linear regression model's parameters (i.e. coefficient) that minimize the sum of squared residual.
- the residuals are differences between observed values of dependent variable and predicted values of dependent variable given w.r.t independent variable
- OLS Algorithm assumes that the errors are normally distributed with zero mean and constant variance and that there is no multicollinearity (high correlation) among the independent variables.
- other method like generalized least square or weighted least square, should be used in case where these assumption are not meet.

Let understand with problem

X	1	2	3	4	5	6	7
y	1.5	3.8	6.7	9.0	11.2	13.6	16.

We will calculate the equation for the best fit line where all the point will be as close as possible by least square method.

X	y	xy	x^2	$\sum x = 28$	$\sum y = 61.8$
1	1.5	1.5	1		
2	3.8	7.6	4	$\sum xy = 314.8$	$\sum x^2 = 140$
3	6.7	20.1	9		
4	9.0	36	16		
5	11.2	56	25		
6	13.6	81.6	36		
7	16	112	49		
Σ	28	314.8	140		
				$y = mx + b$	

$$M = \frac{n \sum xy - \sum x \sum y}{n \sum x^2 - (\sum x)^2} = \frac{7(314.8) - (28)(61.8)}{7(140) - (28)^2}$$

$$m = \frac{473.2}{196} = 2.4142857.$$

$$b = \frac{\sum y - m \sum x}{n} = \frac{61.8 - 2.4142857(28)}{7}$$

$$b = -0.828571.$$

to get linear equation we should plug value in

$$y = mx + b.$$

$$y = 2.41x * \quad y = 2.41x - 0.83.$$

I'll put it for 2

$$y = 2.41(2) - 0.83 = 3.99$$

$$y = 2.41(5) - 0.83 = 11.22$$

$$y = 2.41(7) - 0.83 = 16.09$$

yact

3.8

11.2

16.

Syntax: statsmodel.api.OLS(y, x).

y - dependent variable x - independent variable.

- Import statsmodel.api as sm
import pandas as pd.

- # reading data from csv
df = pd.read_csv('train.csv')

- # defining the variables

x = df['x'].tolist()

y = df['y'].tolist()

adding the constant term α

$x = sm.add_constant(x)$

performing regression and fitting model.

result = sm.OLS(y, x).fit()

print summary table,

print(result.summary())

What is Ridge Regression (L2 Regularization)

Ridge regression is model tuning method, that is used to analyse any data that suffers from multicollinearity. This method performs L2 regularization when issue of multicollinearity occurs, least squares are unbiased, and variance are large. This result in predicted values being far away from actual values.

Regularization:- it is technique used to calibrate machine learning model to minimize adjusted loss function and avoid overfitting and underfitting.

There are three types of regularization techniques:

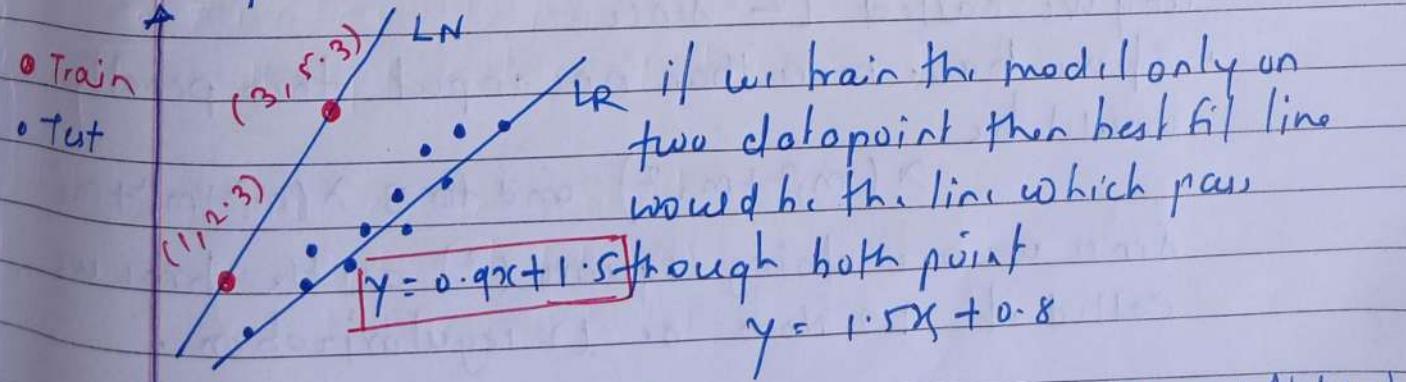
- 1) Ridge Regression (L2 regularization)
- 2) Lasso Regression (L1 regularization)
- 3) Elastic Net (Combo of Ridge and Lasso)

Ridge Regression:-

Overfitting:- means your train acc is too high but test accuracy is very low.

$$y = mx + b$$

m or slope which is coefficient of x is didn't change in y , so to reduce overfitting means to reduce slope.



Now we have to convey our model to choose LR and Not LN

$$L = \sum_{i=1}^n (\hat{y}_i - y_i)^2 + \lambda(m^2) \quad \text{--- (1)}$$

Now we add penalty before we were only reducing which will add help this term which is nothing but Error / Residual for the first term λ = hyperparameter m = slope.

Now we will calculate loss for both line for $\lambda = 1$ with eqn ①

Since line is passing through both points perfectly that's why 1st term will be zero

$$\begin{aligned} L &= 0 + 1(1.5)^2 \\ &= 2.25 \end{aligned}$$

$$\begin{aligned} &(2.3 - 0.9 - 1.5)^2 + \\ &(5.3 - 2.7 - 1.5)^2 + (0.9)^2 \\ &= (0.1)^2 + (1.1)^2 + (0.9)^2 \end{aligned}$$

$$= \boxed{2.03}$$

here we getting significant reduction in loss for the new line

As our model can see this change it will select 2nd model although it will give bad accuracy on training since Variance is significantly reduced although bias increased.

why we called L^2 since

- if we have more than one input then penalty would be

$$\lambda(m_1^2 + m_2^2) \text{ and for } 3 \lambda(m_1^2 + m_2^2 + m_3^2)$$

since we are doing square (^2) all the terms we called it L2 Norm or L2 regularization.

Lasso Regression (L1 Regularization)

- this is also help to reduce overfitting.
- In Ridge regression we have seen.

$$L = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda \|w\|^2$$

MSE penalty term
 $(w_1^2 + w_2^2 + \dots + w_n^2)$
Weight in MLR

Lasso is just another variation of Ridge

$$\begin{aligned}
 L &= \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda \|w\| \\
 &\quad \text{OR} \\
 &= \text{MSE} + \lambda [|w_1| + |w_2| + |w_3| + \dots + |w_n|]
 \end{aligned}$$

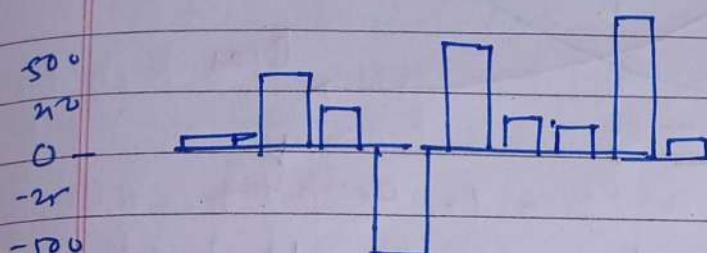
- In ridge regression for any value of λ there were always some value for coefficient of input feature but in Lasso if you continuously increase value of λ at certain point you will zero value for some of coefficient which are not important so here we unknowingly doing feature selection and it is advantage of Lasso.
- so when you are working on high dimensional data and sum. feature are not imp we should prefer Lasso over Ridge.

there are some keypoints need to discuss about Lasso.

1) How coefficients are affected?

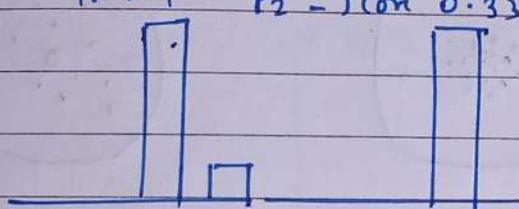
$$\lambda = 0, r_2\text{-score} = 0.44$$

$$\lambda = 0.1, r_2\text{-score} = 0.43$$



$x_1 \ x_2 \ x_3 \ x_4 \ x_5 \ x_6 \ x_7 \ x_8 \ x_9$

$$\lambda = 1, r_2\text{-score} = 0.33$$



$$\lambda = 10, r_2\text{-score} = 0.01$$

{ All coefficients will be zero }

2. Higher coefficients are affected more

- generally as you increased λ coefficient value will be decreased gradually toward zero
- usually higher coefficients are affected first/rapidly
- for it we should be cautious for selecting optimum value to do proper feature selection

3. Impact on Bias and Variance

- As we know if λ increases, overfitting decreases (\downarrow) which lead to increase in bias (\uparrow)

High Bias



Underfitting

High Variance



Overfitting

fig a.

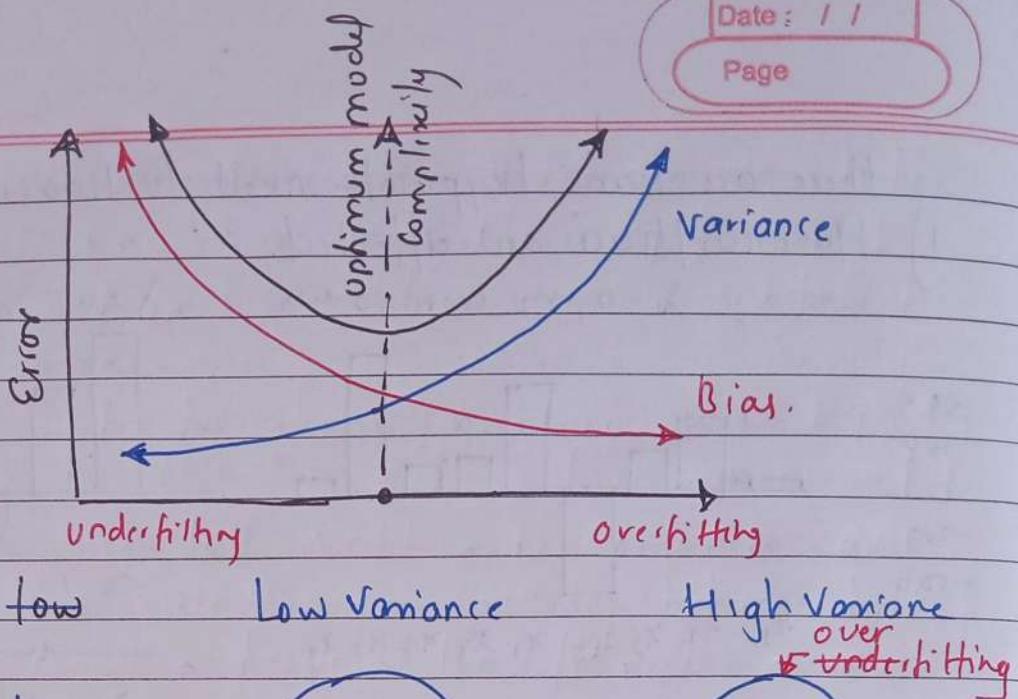
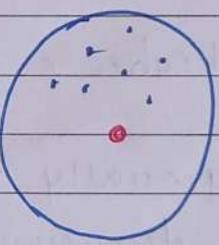
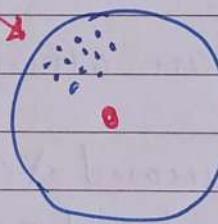


fig (b)

low bias



High bias

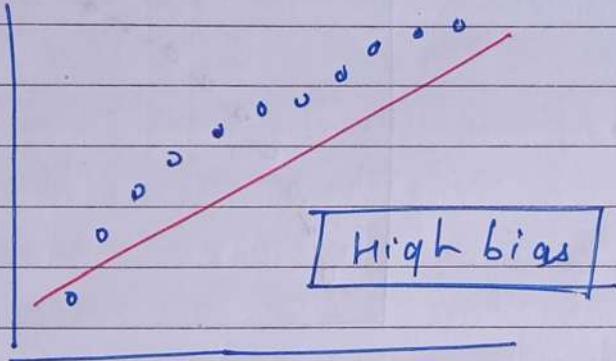


Bias - Variance Tradeoff

- it is important to understand prediction errors (bias and variance) when it comes to accuracy in any ml Algorithm
- There is tradeoff between model's ability to minimize bias and variance which is referred to as best solution for selecting a value of regularization constant
- proper understanding of these errors would help to avoid the overfitting and underfitting of a dataset while training algorithm.

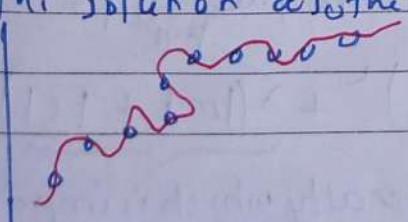
Bias: Bias is known as difference between prediction values by ML model and correct values. Being high in biasing will give large error in training as well as testing. & that's why it is always recommended that algorithm should always be low biased to avoid underfitting.

- if high bias data is predicted in straight line form, thus not fitting accurately in the data in the dataset just fitting called as underfitting.
- this happens the hypothesis is too simple or linear in nature



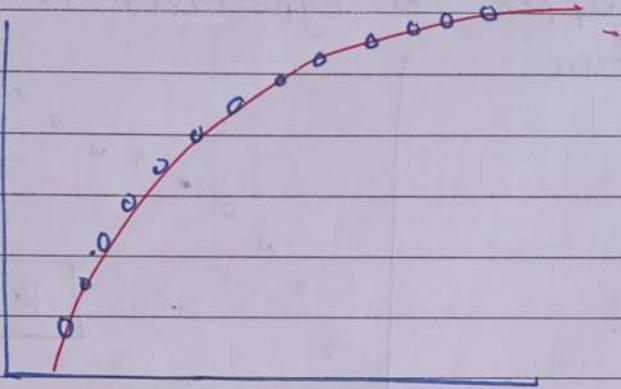
Variance

- The variability of model prediction for given data point which tells us spread of our data is called Variance of model.
- the model with high variance has a very complex fit to training data and thus not able to fit or the test data or data hasn't seen as result such model works very well on training data but has high error rates on test data.
- when model is high on variance it is said to be overfitting of data. Overfitting is fitting the training set accurately via complex curve and high order hypothesis but is not the solution as the error with unseen data is high.



high variance data look like follows

- Bias-Variance tradeoff:- if the algorithm is too simple (hypothesis with linear eqⁿ) then it may be on high bias and low variance and thus it is error prone.
- if error fit too complex (hypothesis with high degree equation) then it may be on high variance and low bias, in the latter condition the new entry will not perform well. there is something between both of these condition known as tradeoff or bias-variance tradeoff



4) Effect of Regularization on loss function loss function:-

$$\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

- it measures how far an estimated value from its true value
- if we are training on different models LR, DT, RF to know which model performs better and which parameters are better loss function is useful.

Elastic Net :- $(L_1 + L_2)$: it is combination of both regularization techniques

$$L_{\text{Reg}} = \frac{\sum (y_i - \hat{y}_i)^2}{n} + \underbrace{\lambda (|m| + l(c))}_{\text{hyperparameter}}$$

penalty which is imposed on norm of

$$L_2 \text{ reg} : \underbrace{\sum_i (y_i - \hat{y}_i)^2}_{n} + \lambda (m^2 + c^2)$$

penalty

Elastic (L1 + L2) :-

$$\sum_i (y_i - \hat{y}_i)^2 + \lambda [(|m| + |c|) + (m^2 + c^2)]$$

adding combination percentage of Lasso & Ridge

$$\sum_i (y_i - \hat{y}_i)^2 + \lambda \left[\underbrace{c (|m| + |c|)}_{\text{Lasso}} + \underbrace{(1-c)(m^2 + c^2)}_{\text{Ridge}} \right]$$

c is between 0 to 1

c = 1 Lasso

c = 0 RIDGE

c = 0.5 = 50% Ridge & 50% Lasso.

Summary :- Ridge is majority is going to focus on regularization but Lasso is going to focus on feature selection as well

- if my job is only feature selection i will go for Lasso and if i want regularization then i will go for Ridge and if i want both then i would go for Elastic Net.

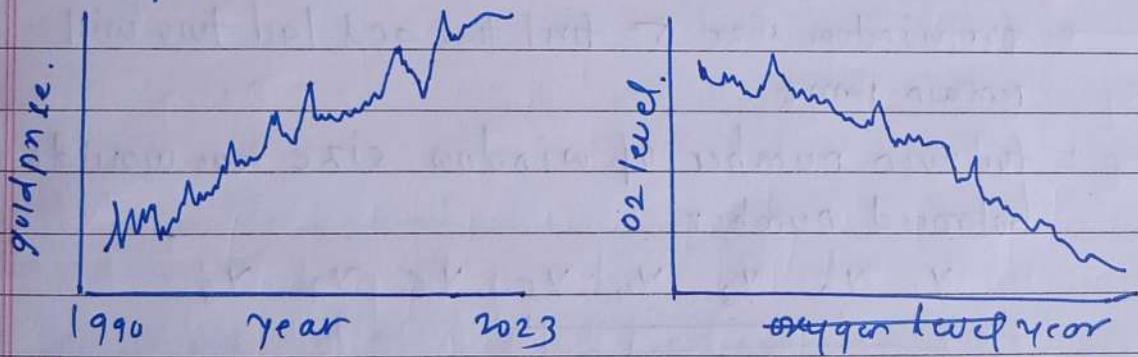
Time Series forecasting.

- time series: - it is data which is index by time.
- the most important part of time series is sequence.

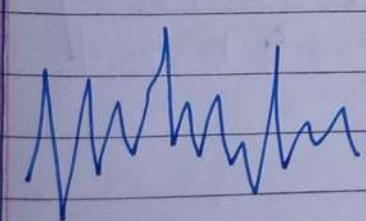
what we do with time series: - we capture the past data & based on that we predict for future this entire procedure is called forecasting

- A) • plotting time series
- component of time series
- forecasting in time series
 - (1) Data based prediction
 - (2) Model based prediction

A) line plot: -



- line plot gives holistic view or overall view.
- line plot gives basic understanding about long term i.e. for long term what happened to my data.
- after time line plot are difficult to understand.



To overcome this we have technique called smoothing.

how to do smoothing :-

moving average:- ex i have rainfall data

Month	M_1	M_2	M_3	M_{120}
data	y_1	y_2	y_3	y_{120}

where M_1 - January

M_2 - february

M_3 - March.

for window size 3

$$y_2 = \frac{y_1 + y_2 + y_3}{3} \quad y_3 = \frac{y_2 + y_3 + y_4}{3} \text{ and so on}$$

- but it would not be same for first and last data point since they will not have one datapoint so their value will remain as it is.
- for window size 5 first two and last two will remain same.
- for even number of window size we wouldn't get balanced number

$$y_1 \ y_2 \ y_3 \ y_4 \ y_5 \ y_6 \ y_7 \ y_8$$

$\underbrace{\qquad\qquad\qquad}_{4} \quad \underbrace{\qquad\qquad\qquad}_{4}$

$$y_3 = \frac{y_2 + y_3 + y_4 + y_5}{4} \quad \{ \text{giving high priority to future data} \}$$

$$y_3 = \frac{y_1 + y_2 + y_3 + y_4}{4} \quad \{ \text{giving high priority to past data} \}$$

Solution is:-

$$= \frac{1}{2} \left[\frac{y_1 + y_2 + y_3 + y_4}{4} \right] + \frac{1}{2} \left[\frac{y_2 + y_3 + y_4 + y_5}{4} \right]$$

• How to figure out best window size?

1. Way of domain knowledge.

Eg. sales of AC is based on season like in summer it sold out rapidly whereas in winter is not seem that like expected.

so here we can assume or try window size 3/4

2. When we don't have domain knowledge another way it by calculating size for all and check where it becomes smooth.

Components of time series:

it generally has three components:

1) Trend - holistic view - Moving average

2) Seasonality - periodic change

3) Ruid - Noise or Error - lesser the better.

Decomposition of T.S [Trend, Seasonality, Noise]

• Trend: - find optimum window size

• To find seasonality & Noise subtract trend from original data.

Let say we have data 2001 to 2010

Month data	M_1	M_2	M_3	M_4	\dots	\dots	\dots	M_{120}
	y_1	y_2	y_3	y_4	\dots	\dots	\dots	y_{120}
Trend	$T_1 = M_1$	$T_2 = \frac{M_1 + M_2 + M_3}{3}$	$T_3 = \frac{M_2 + M_3 + M_4}{3}$	\dots	\dots	\dots	\dots	\dots

$$\text{Seasonality} : S_1 = \text{Jan} = \frac{T_1 + T_{13} + T_{25} + \dots + T_{109}}{10}$$

$$S_2 = \text{Feb} = \frac{T_2 + T_{14} + T_{26} + \dots + T_{110}}{10}$$

$$\text{Noise} : N_1 = y_1 - T_1 - S_1, N_2 = y_2 - T_2 - S_2, \dots, N_{120} = y_{120} - T_{120} - S_{120}$$

forecasting in time series.

Data based forecasting :-

- ① Simple Exponential :- it is applicable to those time series where there is no trend no seasonality.
- ② Double Exponential :- it is only performed good for those time series where it only have trend but not seasonality.
- ③ Holt winter's model - (Triple Exponential) it is performed well where there is both trend as well as seasonality (when you don't know whether time series have trend or seasonality then holt winter is best)

→ Simple Exponential :-

Day	D ₁	D ₂	D ₃	D ₁₉₉	D ₂₀₀	D ₂₀₁
Temp	t ₁	t ₂	t ₃	t ₁₉₉	t ₂₀₀	Today → (predict)

what will be the easiest way to predict for tomorrow's temp when we have data upto today?

i) taking Avg :- $\frac{t_1 + t_2 + t_3 + \dots + t_{200}}{200} = n$.

but this not good method since temp varies a lot throughout year summer to winter

ii) another way is whatever is today's temp it will same for tomorrow.

$$t_{\text{day 200}} = 25.7^\circ$$

$$t_{\text{day 201}} = 25.7^\circ$$

- Simple Exponential:- instead of taking all 200 data into consideration, we will take recent data for prediction.
- Simple Exponential try to calculate the component in Local Average.

Day	D_1	D_2	D_3	D_4	.	.	.	D_{199}	D_{200}
temp	t_1	t_2	t_3	t_4				t_{199}	t_{200}
	L_1	L_2	L_3	L_4				L_{199}	L_{200}

- Local Average at Day 4 i.e. L_4 will be combination of two component yesterday's Local Average and today temp

$$L_3 \xrightarrow{L_4} t_4 = L_4 = \frac{1}{2} L_3 + \frac{1}{2} t_4$$

Note: $L \rightarrow t_1$ since there is no previous record.

Simple Exponential is applicable second point onwards.

Mathematical Understanding :-

$$L_{200} = \frac{1}{2} L_{199} + \frac{1}{2} t_{200} \quad \textcircled{1}$$

$$L_{199} = \frac{1}{2} L_{198} + \frac{1}{2} t_{199} \quad \textcircled{2}$$

$$L_{200} = \frac{1}{2} \left[\frac{1}{2} L_{198} + \frac{1}{2} t_{199} \right] + \frac{1}{2} t_{200} \dots \text{2 in } \textcircled{1}$$

$$= \frac{1}{4} L_{198} + \frac{1}{4} t_{199} + \frac{1}{2} t_{200} \quad \textcircled{3}$$

$$L_{198} = \frac{1}{2} L_{197} + \frac{1}{2} t_{198} \quad \textcircled{4}$$

put $\textcircled{4}$ in $\textcircled{3}$

$$\frac{1}{4} \left[\frac{1}{2} L_{197} + \frac{1}{2} t_{198} \right] + \frac{1}{4} t_{199} + \frac{1}{2} t_{200}$$

$$\frac{1}{8} L_{197} + \frac{1}{8} t_{198} + \frac{1}{4} t_{199} + \frac{1}{2} t_{200}$$

If I rearrange this, we will get

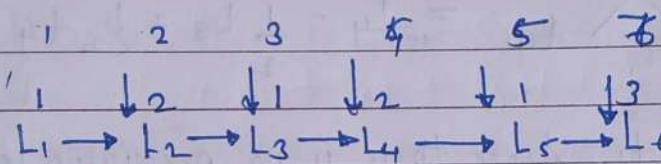
$$L_{200} = \frac{1}{2} t_{200} + \frac{1}{4} t_{199} + \frac{1}{8} t_{198} + \frac{1}{16} t_{197} \dots$$

\uparrow \uparrow \uparrow \uparrow
 $(\frac{1}{2})^1$ $(\frac{1}{2})^2$ $(\frac{1}{2})^3$ $(\frac{1}{2})^4$

from current datapoint weightage of coefficient reduce exponentially whereas in simple average every coefficient were getting equal average.

Let understand with short Example.

Day
Score



$$L_1 = 1, L_2 = \frac{L_1 + S_2}{2} = \frac{1+2}{2} = 1.5$$

$$\underline{\underline{L_3 = \frac{L_2 + S_3}{2} = \frac{1.5+1}{2}} \quad L_4 = \frac{L_3 + S_4}{2} = \frac{1.25+2}{2} = \frac{3.25}{2}}$$

$$= \frac{2.5}{2} = 1.25 \quad = 1.625$$

$$\underline{\underline{L_5 = 1.3125, L_6 = 2.1562}}$$

All future prediction going to be last available local Average.

Double Exponential :- it contain trend but no seasonality.

We can't always rely on local Average because for ex. if chocolate price increase by Rs 1 each day.

Day	1	2	3	4	5	6
Price	1	2	3	4	5	6

L_6

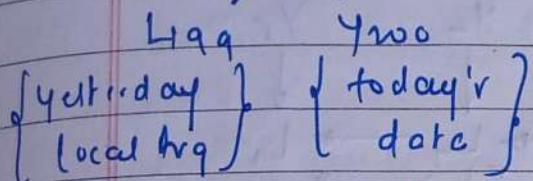
L_6 = around 4.2 which is not true.

Hence we need consider a trend as well, that's why we need both trend as well as Local Average.

for local Avg

but local avg is also

expected to change by some amount.



there is little modification in local Average

$$L_{200} = \frac{1}{2} L_{199} + \frac{1}{2} Y_{200}$$

$$L_{200} = \frac{1}{2} (L_{199} + T_{199}) + \frac{1}{2} Y_{200}$$

↑
trend is added.

$$T_{200} = \frac{1}{2} T_{199} + \underbrace{\frac{1}{2} (L_{200} - L_{199})}_{\text{change in local Average from yesterday to today}}$$

↑
(today's)
trend

Raw data	y_1	y_2	y_3	y_4	\dots	y_{200}
Local Avg	L_1	L_2	L_3	L_4	\dots	L_{200}
Trend	T_1	T_2	T_3	T_4	\dots	T_{200}

so from raw data we will calculate Local Avg & trend.

$$L_1 = 0 \quad L_2 = \frac{1}{2} (L_1 + T_1) + \frac{1}{2} (y_2)$$

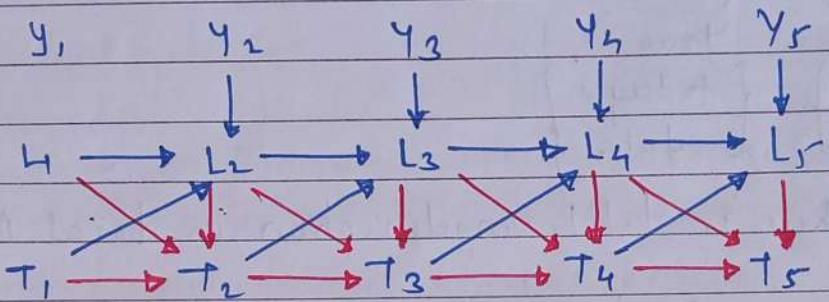
$$T_1 = 0$$

$$T_2 = \frac{1}{2} (T_1) + \frac{1}{2} (L_2 - L_1).$$

Let us understand with Example, we will solve problem, both Simple Exponential & Double Exponential.

Day	1	2	3	4	5	$\frac{1+2}{2} = 1.5$
Price (Y)	1	2	3	4	5	
L(S)	1	1.5	2.25	3.125	4.0625	$\frac{1.5+3}{2} = 4.5 \div 2 = 2.25$
L(D)	1	1.5	2.375	3.46875	4.6484375	

Now we will calculate with Double Exponential.



$$L_2 = \frac{1}{2}(L_1 + T_1) + \frac{1}{2}(Y_2)$$

$$T_2 = \frac{1}{2}(T_1) + \frac{1}{2}(L_2 - L_1)$$

$$L_2 = \frac{1}{2}(1+0) + \frac{1}{2}(2) = \frac{1}{2} + 1 = 1.5$$

$$T_2 = \frac{1}{2}(0) + \frac{1}{2}(1.5-1) = 0 + \frac{0.5}{2} = 0.25$$

$$L_3 = \frac{1}{2}(L_2 + T_2) + \frac{1}{2}(Y_3)$$

$$\Rightarrow \frac{1}{2}(1.5 + 0.25) + \frac{1}{2}(3) = \frac{1}{2}(1.75) + 1.5$$

$$= 0.875 + 1.5$$

$$= 2.375$$

$$T_3 = \frac{1}{2}T_2 + \frac{1}{2}(L_3 - L_2)$$

$$= \frac{1}{2}(0.25) + \frac{1}{2}(2.375 - 1.5)$$

$$= 0.125 + 0.4375$$

$$= 0.5625$$

$$\begin{aligned}
 L_4 &= \frac{1}{2}(L_3 + T_3) + \frac{1}{2}(\gamma_4) \\
 &= \frac{1}{2}(2.375 + 0.5625) + \frac{1}{2}(5) \\
 &= \frac{2.9375}{2} + 2 \\
 &= 1.46875 + 2 = 3.46875
 \end{aligned}$$

$$\begin{aligned}
 T_4 &= \frac{1}{2}(T_3) + \frac{1}{2}(L_4 - L_3) \\
 &= \frac{1}{2}(0.5625) + \frac{1}{2}(3.46875 - 2.375) \\
 &= 0.28125 + 0.5(1.09375) \\
 &= 0.28125 + 0.546875 \\
 &= 0.828125
 \end{aligned}$$

$$\begin{aligned}
 L_5 &= \frac{1}{2}(L_4 + T_4) + \frac{1}{2}(\gamma_5) \\
 &= \frac{1}{2}(3.46875 + 0.828125) + \frac{1}{2}(5)
 \end{aligned}$$

$$= 2.1484375 + 2.5$$

$$= 4.6484375$$

$$\begin{aligned}
 T_5 &= \frac{1}{2}T_4 + \frac{1}{2}(L_5 - L_3) \\
 &= \frac{1}{2}(0.828125) + \frac{1}{2}(4.6484375 - 2.375) \\
 &= 0.4140625 + 0.5(1.1797) \\
 &= 0.4140625 + 0.589875 \\
 &= 1.0039375
 \end{aligned}$$

$$\begin{array}{ccccccccc}
 T_1 & T_2 & T_3 & T_4 & T_5 & & \} & \text{Trend is increased.} \\
 0 & 0.25 & 0.56 & 0.82 & 1.0039 & & &
 \end{array}$$

Actual Predicted

4.0625

5

Simpl. Exponential

4.6484

5

Double Exponential

4.6484

5

Holt's winter model :- which contain both trend & seasonality.

- Here along with local Average we will find out trend as well as seasonality.

time series data

y_1	y_2	y_3	y_{200}
L_1	L_2	L_3	L_{200}
trend T_1	T_2	T_3	T_{200}
seasonality s_1	s_2	s_3	s_{200}

for 200 Months \rightarrow 16 years - 7 months

for seasonality - periodic h.m. of 12 months

Let first convert it into double exponential problem.
by subtracting seasonality from time series

$$y_1 - s_1, y_2 - s_2, y_3 - s_3 \dots \dots \dots y_{200} - s_{200}$$

$$L_{200} = \frac{1}{2}(L_{199} + T_{199}) + \frac{1}{2}(y_{200} - s_{200}) \quad \text{--- (1)}$$

{ Note Since $y_{200} = y_{200} - s_{200}$ }

$$T_{200} = \frac{1}{2}(T_{199}) + \frac{1}{2}(L_{200} - L_{199}) \quad \text{--- (2)}$$

$$s_{200} = \frac{1}{2}(s_{200} - 12) + \frac{1}{2}(s_{200})$$

From (1) we can get to know to get
 L_{200} we subtract s_{200} from y_{200} then
now s_{200} we can do $y_{200} - L_{200}$.

$$s_{200} = \frac{1}{2}(s_{200} - 12) + \frac{1}{2}(y_{200} - L_{200}) \quad \text{--- (3)}$$

If you are carefully observe then you find there is loop between ③ equations.

for L_{200} we need s_{200} for s_{200} we need L_{200}

Since seasonality is periodic

$$L_{200} = s_{188}$$

in eq^① $s_{200} = s_{188}$

$$L_{200} = \frac{1}{2} (L_{99} + T_{199}) + \frac{1}{2} (Y_{200} - s_{188})$$

$$T_{200} = \frac{1}{2} (T_{199}) + \frac{1}{2} (L_{200} - L_{199})$$

$$s_{200} = \frac{1}{2} (s_{200} - l_2) + \frac{1}{2} (Y_{200} - L_{200})$$

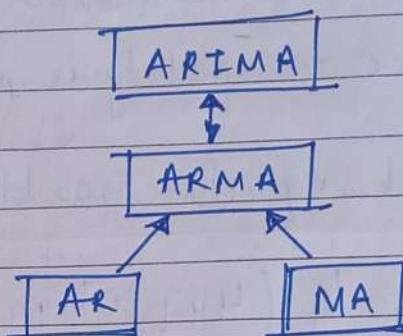
Summary table.

Model	when we can apply	what should be future prediction	hyperparameters
Simple	No trend	last possible local Avg = 1	
Exponential	No seasonality.	all future prediction	
Double Exponential	only trend	last possible local Avg & 2	
	No seasonality	changed trend	
Triple Exponential	both trend and local Average + trend + 3.	seasonality.	

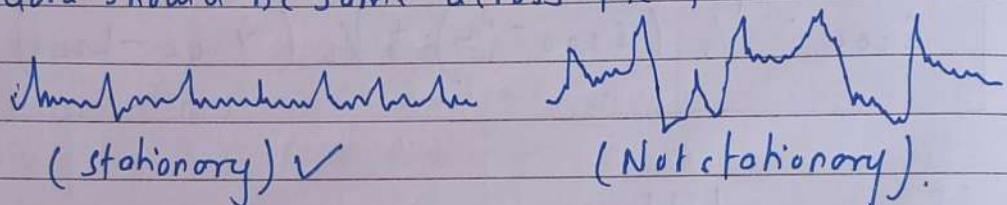
Model Based forecasting :-

Basically there are three models.

- 1) AR
- 2) MA
- 3) ARMA
- 4) ARIMA



Here data must be stationary. It means distribution of data should be same across the time.



If we have stationary data then go for model based otherwise databased, since there is no any condition

AR (Auto Regressive) : AR is completely depend upon past data and nothing else for future prediction

MA: Moving Average : future data is only depend only external outsiders factors.

ARMA:- future data only going to deal with both past data data and external factors.

AR Model :-

$$\text{In regression } y = \alpha x + \beta + \epsilon$$

$$\text{future } y_t = \alpha y_{t-1} + \beta + \epsilon$$

↑ ↑ ↑ ↑

past y Beta Error

Here future data is only depend on past data
if it would depend on past two data

$$Y_t = \alpha_1 Y_{t-1} + \alpha_2 Y_{t-2} + \beta + \varepsilon \quad AR(2)$$

$$Y_t = \alpha_1 Y_{t-1} + \alpha_2 Y_{t-2} + \alpha_3 Y_{t-3} + \beta + \varepsilon \dots A(3)$$

How to select Best possible model. $AR1, AR2, AR3$
~~AR1~~ based on test data accuracy whichever will perform good will finalize that.

MA : Moving Average Here my future data is only depend upon external factors.

Ex : rate of some production which depends upon seasons. (grain or fruit) due to natural calamities.

data	y_1	y_2	y_3	-	-	y_t
External factor	ε_1	ε_2	ε_3	-	-	ε_t

$$y_t = \alpha \varepsilon_{t-1} + \beta + \varepsilon_t \dots MA(1)$$

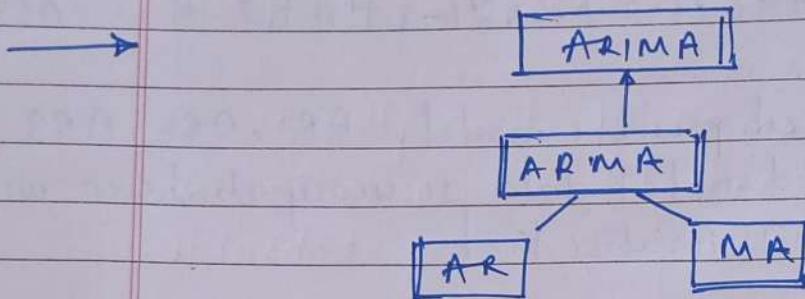
this year data | last year data \ constant

$$y_t = \alpha_1 \varepsilon_{t-1} + \alpha_2 \varepsilon_{t-2} + \beta + \varepsilon_t \dots MA(2)$$

Model based forecasting.

- ① AR
- ② MA
- ③ ARMA (AR + MA)

④ ARIMA (which is combination of ARMA itself)



AR → future prediction are completely depend on past data, it is just copy of regression model.

$$y = \alpha x + \beta + \epsilon$$

target variable prediction error.

$$y_t = \alpha y_{t-1} + \beta + \epsilon \quad \text{AR(1)}$$

$\underbrace{\alpha y_{t-1}}_{\text{past data}}$ $\underbrace{\beta}_{\text{it is depend on one past data}}$

$$y_t = \alpha_1 y_{t-1} + \alpha_2 y_{t-2} + \alpha_3 y_{t-3} + \beta + \epsilon$$

\uparrow tomorrow \uparrow today \uparrow yesterday \uparrow day before yesterday

if it is similar to SLR

$$\text{SLR: } y = \alpha x + \beta + \epsilon$$

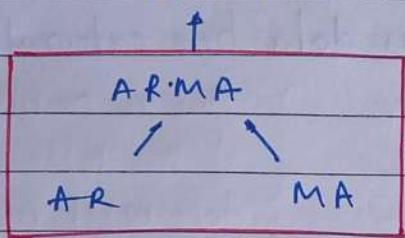
$$\text{AR}(1) = \alpha y_{t-1} + \beta + \epsilon$$

$$\text{2 features: } y = \alpha_1 x_1 + \alpha_2 x_2 + \beta + \epsilon$$

$$\text{AR}(2) = \alpha_1 y_{t-1} + \alpha_2 y_{t-2} + \beta + \epsilon$$

- one of the necessary condition for apply these data is data should be stationary i.e. it means distribution of data should be same across all the data point.
- Since among all four model ARMA, AR and MA are very critical for that we must sure that data should be stationary.

ARIMA



Since they are critical
must sure for stationary
data.

MA → Moving Average :- here future prediction are completely depend upon external factors.

production of rice	y_1	y_2	y_3	y_4	y_5	y_6	y_7
Error due to external factor	ϵ_1	ϵ_2	ϵ_3	ϵ_4	ϵ_5	ϵ_6	ϵ_7

- when we say our data is only depending on last time point for example for 8th year

$$y_8 = \alpha \epsilon_7 + \beta + \epsilon_8$$

so general formula is

$$y_t = \alpha y_{t-1} + \beta + \epsilon_t \quad \text{Moving Avg (MA)}$$

depend on External Error factors.

$$y_t = \alpha y_{t-1} + \beta + \epsilon_t \quad \text{Auto Regressiv (AR)}$$

depend on past data.

for MA(2) and MA(3) equation would be.

$$MA(2): Y_t = \alpha_1 \epsilon_{t-1} + \alpha_2 \epsilon_{t-2} + \beta + \epsilon$$

$$MA(3): Y_t = \alpha_1 \epsilon_{t-1} + \alpha_2 \epsilon_{t-2} + \alpha_3 \epsilon_{t-3} + \beta + \epsilon.$$

- it is exactly similar to AR model only difference is instead of previous data here external factors are introduced.

|| ARMA ||: it is combination of AR and MA model.

ARMA (1, 1)		
AR (1)	+	MA (1)

$$ARMA(p, q) = AR(p) + MA(q)$$

$$Y_t = \alpha_1 Y_{t-1} + \underline{\beta_1 \epsilon_{t-1}} + \beta_0 + \epsilon_t.$$

if $\beta_1 \epsilon_{t-1} = 0$ then it will be AR (1)

$$Y_t = \underline{\alpha_1 t_{t-1}} + \beta_1 \epsilon_{t-1} + \beta_0 + \epsilon_t.$$

if $\alpha_1 t_{t-1} = 0$ then it will be MA (1)

- that's why you don't need to put MA and AR model separately.
- if you think your model mostly inclined to MA then we can put AR component zero if model inclined to AR then we can put MA component zero.

- thing we need to take care for these three models.
- i) stationary ... if it is stationary we gonna directly used ARMA instead of using AR and MA Model separately.
- ii) ARMA model has two hyperparameter (P, q)
 P corresponds to AR and q corresponds to MA

how to choose best value of P and q .

by creating multiple model in range of (0 to 5)
for each P and q

$$P = \{0, 1, 2, 3, 4, 5\} = 6 \quad 6 \times 6 = 36$$

$$q = \{0, 1, 2, 3, 4, 5\} = 6$$

since 0,0 is not possible pair we have $36 - 1 = 35$

ARIMA : ARIMA doesn't care about data is stationary or Not.

this not a new model but extension of ARMA

ARIMA - Auto Regressive integrated moving average.

- if transformed data from original data follows ARMA model then original data will must follow ARIMA.

$$\begin{array}{ccccccc} y_1 & y_2 & y_3 & y_4 & \dots & y_n \\ \downarrow & \downarrow & & & & & \downarrow \\ z_1 & z_2 & z_3 & z_4 & \dots & z_n \end{array}$$

if this $\{z_1, z_2, \dots, z_n\}$ follows ARMA then

(y_1, y_2, \dots, y_n) will must follow ARIMA.

for ex.

Bank balance of Every month

$y_1 \ y_2 \ y_3 \ y_4 \dots \ y_t \leftarrow (\text{today})$

{ here money is added every month not due to new
added but due to interest}

what will interest added each month

$\downarrow \quad \downarrow \quad \downarrow \quad \downarrow \quad \downarrow$
 $Z_1 \ Z_2 \ Z_3 \ Z_4 \ Z_t$

where $Z_1 = y_2 - y_1$ } this is newly constructed data
 $Z_2 = y_3 - y_2$ } if this follow ARMA then the parat
 $Z_3 = y_4 - y_3$ } data will follow ARIMA

ARMA has two component (P, q) whereas

ARIMA has three component (P, q, d)

¹ degree of subtraction

d - for subtraction of consecutive numbers.

$d = 1$

$d = 2$

$d = 3$

$$Z_1 = y_2 - y_1$$

$$Z_1 = y_3 - y_1$$

$$Z_1 = y_4 - y_1$$

$$Z_1 = y_3 - y_2$$

$$Z_2 = y_4 - y_2$$

$$Z_1 = y_5 - y_2$$

$$\vdots \quad \vdots \quad \vdots$$

$$\vdots \quad \vdots \quad \vdots$$

$$\vdots \quad \vdots \quad \vdots$$

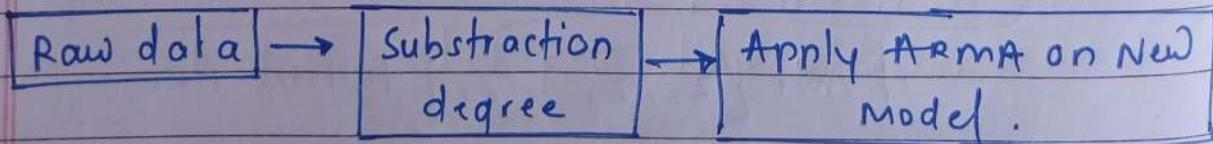
$$Z_{100} \ y_{101} \ y_{100}$$

$$Z_{100} = y_{102} - y_{100}$$

$$Z_{100} = y_{103} - y_{100}$$

what is lifecycle of ARIMA

Transformation



Summary table.

Model No of parameter Conditions depends on.

AR 1 stationary past

MA 2 stationary External factor

ARMA 2 stationary past +
External factor

ARIMA 3 No need - since I am
transforming the data.

forecasting

Data Based

there is no previous
condition

Simple Exponential Double Exponential Holt winter

No trend trend ✓ both
No seasonality. but No seasonality. trend & seasonality.

Model Based

AR MA ARMA ARIMA
[stationary data]

{ when we are not sure }

directly go for ARIMA

{ whenever we don't
know which to
use we can used
directly Holt winter }

Evaluation of time series Model.

there are two matrix we gonna used.

- ① Mean Square Error
- ② MAPE score.

Mean Square Error.

let say i have so many future data.

Actual	Prediction
y_t	\hat{y}_t
y_{t+1}	\hat{y}_{t+1}
y_{t+2}	\hat{y}_{t+2}

$$MSE = \frac{(y_t - \hat{y}_t)^2 + (y_{t+1} - \hat{y}_{t+1})^2 + (y_{t+2} - \hat{y}_{t+2})^2}{n}$$

- ② MAPE Score. (Mean Absolute Percentage Error)

let compare 2 scenarios.

	S1	S2
Actual	3	1000
Predicted	2	999

which is more accurate?

{ if you try MSE for both it would be same }
 it means, sometimes we are not interested in amount of error but percentage of error

formula for MAPE scores is

$$\frac{|y_A - y_P|}{y_A} \quad \text{--- this is value of percentage}$$

$$\text{Final calc. : } \frac{|3-2|}{3} = \frac{1}{3} = 33.3\%$$

$$\text{Second case: } \left| \frac{1000 - 999}{100} \right| = 0.001 \\ 0.01\%$$

Since we are not only concerned about Absolute error always but absolute percentage error also.

Exercises :- calculate the MSE and MAPE score

Actual	Predicted
5	10
20	15
30	25

$$\text{MSE} = \frac{(5-10)^2 + (20-15)^2 + (30-25)^2}{3} \\ = \frac{(-5)^2 + (5)^2 + (5)^2}{3} = 25$$

MAPE score :

$$\left| \frac{y_A - y_P}{y_A} \right| = \left| \frac{5 - 10}{5} \right| = -5/5 = -1 = 100\%$$

$$\therefore \quad = \left| \frac{20 - 15}{20} \right| = 25\%$$

$$\therefore \quad = \left| \frac{30 - 25}{25} \right| = 16.6\%$$

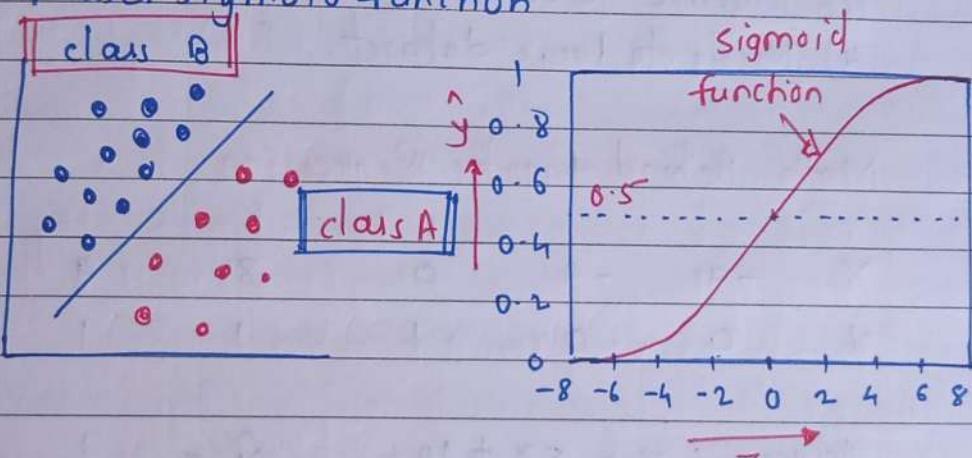
$$\text{Mean abs error} = \frac{100 + 25 + 16.6}{3} = \underline{\underline{47.2}}$$

Classification Algorithm

- 1) Logistic Regression
- 2) SVM
- 3) KNN
- 4) Decision Tree
- 5) Random forest
- .

1) Logistic Regression:-

1. it is supervised learning model.
2. it is classification model and best for binary classification.
3. it uses sigmoid function



$$\text{Sigmoid function} = \hat{y} = \frac{1}{1 + e^{-z}}$$

where $z = w_1x_1 + w_2x_2 + b$.

$(w_1x_1 + w_2x_2 + b)$ - eqn of line

slope $m \rightarrow$ weight w

intercept $c \rightarrow$ bias b

• \hat{y} = probability that ($y=1$)

$\hat{y} = p(y=1|x) \dots \{ \text{probability of } y \text{ being 1 for given value of } x \}$

x = input features

w = weights (it will be in the format)

{ number of weight equal to number of feature in the dataset }

b = bias

$\hat{y} = \sigma(z)$

Advantages: 1) Easy to implement

2) perform well on data with linear relationship
3). less prone to overfitting for low dimensional dataset.

Disadvantages :- 1) High dimensional dataset causes overfitting.

- 2) difficult to capture complex relationship in dataset
- 3) sensitive to outlier
- 4) Needs large dataset.

Moth Behind Logistic Regression.

X	-9	-8	0	8	9
Y	0	0	1	1	1

Assume $Z = 5x + 10$ $\hat{y} = \frac{1}{1 + e^{-Z}}$

$x = -9$	$x = -8$	$x = 0$	$x = 8$	$x = 9$
----------	----------	---------	---------	---------

$$\begin{aligned} Z &= 5(-9) + 10 & Z &= 5(-8) + 10 & Z &= 5(0) + 10 & Z &= 5(8) + 10 & Z &= 5(9) + 10 \\ &= -35 & & = -30 & & = 10 & & = 50 & & = 55 \end{aligned}$$

$$\begin{aligned} \hat{y} &= \frac{1}{1 + e^{-35}} & \hat{y} &= \frac{1}{1 + e^{-30}} & \hat{y} &= \frac{1}{1 + e^{10}} & \hat{y} &= \frac{1}{1 + e^{50}} & \hat{y} &= \frac{1}{1 + e^{55}} \end{aligned}$$

$$\begin{aligned} \hat{y} &= 0 & \hat{y} &= 0 & \hat{y} &= 1 & \hat{y} &= 1 & \hat{y} &= 1. \end{aligned}$$

Inference : if z value is large positive number,

$$\hat{y} = \frac{1}{1+e^{-z}} \approx \hat{y} = 1.$$

if z is large negative number,

$$\hat{y} = \frac{1}{1 + (\text{large positive number})}$$

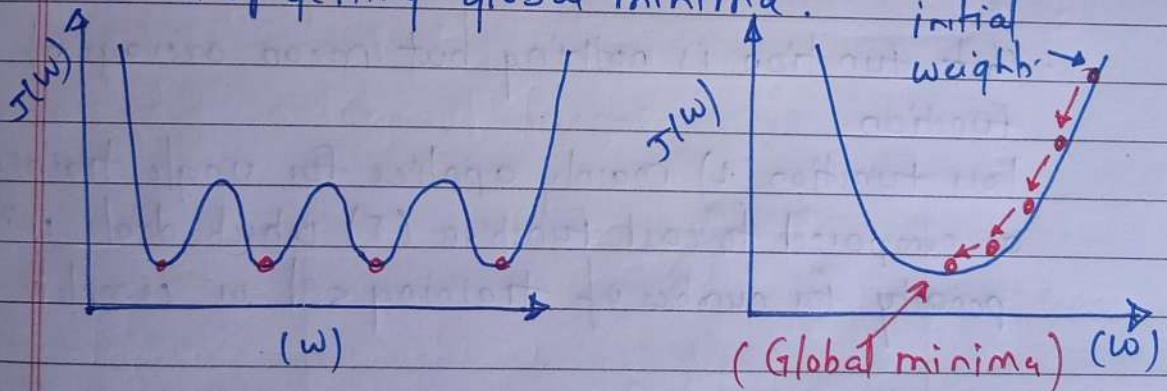
$$\hat{y} = 0.$$

loss function & cost function for Logistic Regression.

- loss function measures how far an estimated value is from true value.

$$\text{loss function for linear regression} = \frac{1}{n} \sum_{i=1}^n (y_i - y_{\text{pred}})^2$$

if we use this function we will get many local minima instead of getting global minima.



- Binary cross entropy loss function (or) log loss.

$$L(y, \hat{y}) = -(y \log \hat{y} + (1-y) \log(1-\hat{y}))$$

Here,

$$y \rightarrow 0 \text{ or } 1$$

$$\text{by } \hat{y} \rightarrow 0 \text{ to } 1 \text{ (probability could be continuous)}$$

when $y = 1$

$$\begin{aligned} L(1, \hat{y}) &= -(1 \log \hat{y} + (1-1) \log(1-\hat{y})) \\ &= -\log \hat{y} \end{aligned}$$

- Since we always want smaller loss function value hence \hat{y} should be very large (from 0 to 1) if it is the $-\log \hat{y}$ will be very large negative number or very small number.

when $y = 0$

$$\begin{aligned} L(1, \hat{y}) &= -(0 \log \hat{y} + (1-0) \log(1-\hat{y})) \\ &= -\log(1-\hat{y}). \end{aligned}$$

- Since we want smaller loss function value, hence \hat{y} should be very small the automatically $(1-\hat{y})$ will be very large thus $-\log(1-\hat{y})$ will be large negative number or very small number.

- Cost function is nothing but mean average of loss function

- Loss function (L) mainly applies for single training set as compared to cost function (J) which deals with a penalty for number of training set or complete batch.

Loss function:

$$L(y, \hat{y}) = -(y \log \hat{y} + (1-y) \log(1-\hat{y})) \quad \text{for single}$$

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m (L(y^{(i)}, \hat{y}^{(i)})) =$$

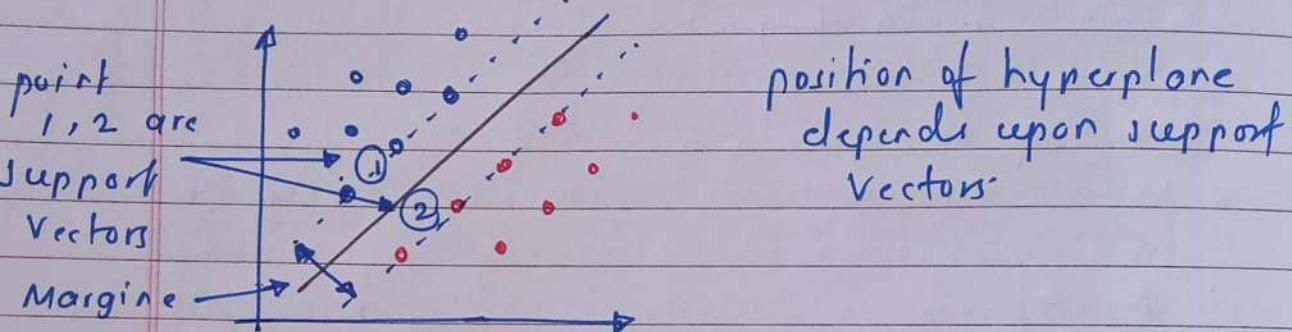
$$-\frac{1}{m} \sum_{i=1}^m (y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log(1-\hat{y}^{(i)}))$$

{'m' denotes number of data points in the }
training set

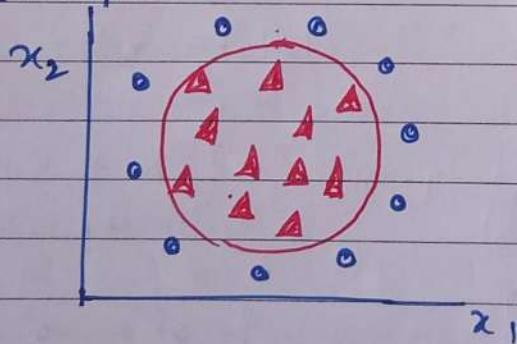
Support Vector Machine (SVM)

Basic about SVM.

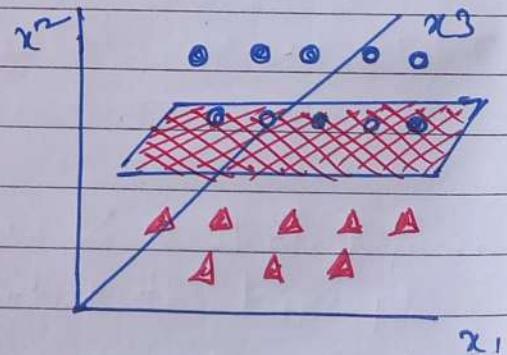
- 1) it is supervised ML model.
- 2) it can be used for both classification as well as regression but it is predominantly used for binary classification.
- 3) Hyperplane.
- 4) Support vectors



- for 2D data it is easy to draw hyperplane but where data point can't be separated by line need to convert into 3D where we can separate the datapoint by hyperplane.
- Example:-



{Here not easy to separate by line}



$2D \rightarrow 3D$
and separate by hyperplane

Hyperplane— Hyperplane is line (in 2D) or plane that separates the data point into two classes

Support Vectors :- these are the datapoints which are nearest to hyperplane if these datapoints change position of hyperplane changes.

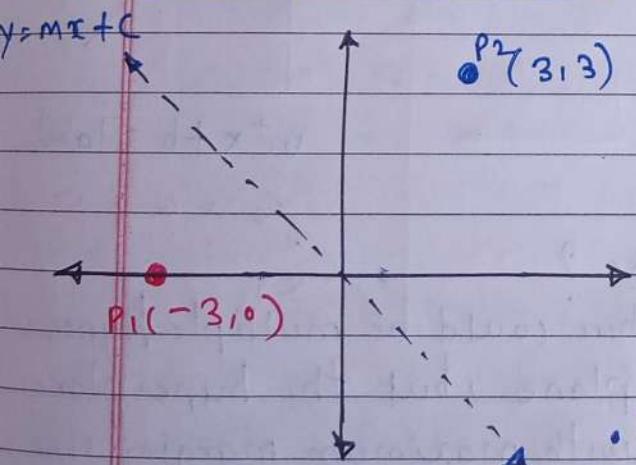
Advantages of SVM

- 1) works fine with smaller dataset
- 2) works fine or efficiently where there is clear margin of separation
- 3) works well with high dimensional data

Disadvantages

- 1) Not suitable for large dataset as training time would take very long.
- 2) Not suitable for noiser (outlier) dataset with overlapping classes.

Math Behind SVM



Let slope and intercept of hyperplane is,
 $m = -1$

$c = 0$ {since passing through origin}

• let parameters of hyperplane sare in w which is nothing but weight
 $w \rightarrow (m, c) = (-1, 0)$

• Let multiply x or P_1 by transpose of w

$$w^T x = [-1] [-3 \ 0] = 3 \text{ (positive)}$$

[Note : why transpose ? \rightarrow for matrix multiplication no of column of 1st Matrix must be equal to no. of rows of second matrix]

- positive value indicates all the points of hyperplane will be positive class.

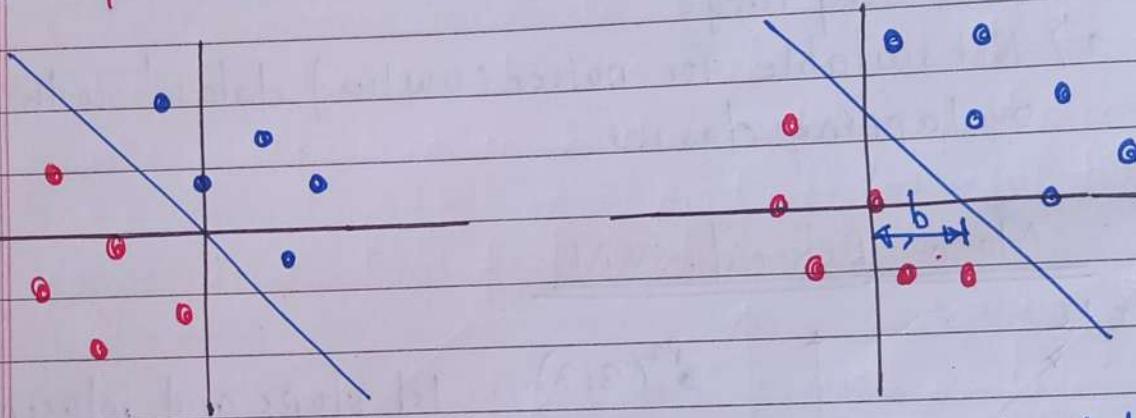
for $P_2(3, 3)$

$$w^T x = \begin{bmatrix} -1 \\ 0 \end{bmatrix} \begin{bmatrix} 3 \\ 3 \end{bmatrix}$$

$$= -3 \text{ (Negative)}$$

Here for all the points which lie on the right side of hyperplane will belong to negative class.

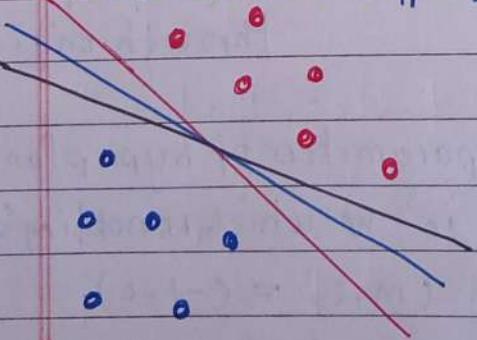
But Not all the time hyperplane will pass through Origin.



$$w^T x = \text{label}$$

$$w^T x + b = \text{label}$$

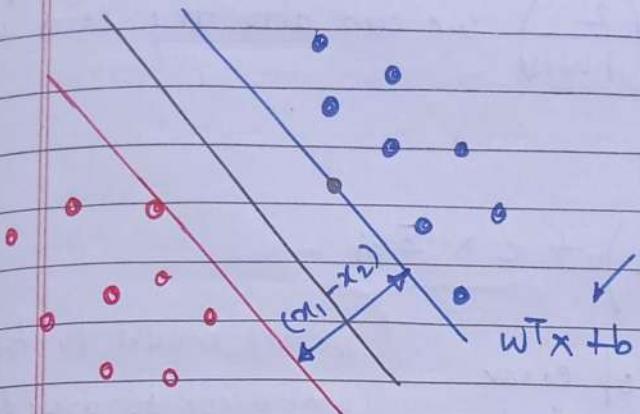
which is best hyperplane?



there could be multiple hyperplane, but the hyperplane with maximum margin size will be the best hyperplane.

→ Optimization for Maximum margin

$$w^T x + b = \text{label}$$



Equation of point or
blue support vector &
its output value any
negative value.

- $w^T x + b = 1 \Rightarrow$ this is equation of point or red support vector and its output value could be any positive value

to get margin let subtract one from another.

$$w^T x_1 + b = 1$$

$$\underline{(-)} \quad w^T x_2 + b = -1$$

$$w^T (x_1 - x_2) = 2$$

$$w^T (x_1 - x_2) = 2$$

divide both sides by $\|w\|$

$$\frac{w^T (x_1 - x_2)}{\|w\|} = \frac{2}{\|w\|}$$

$$(x_1 - x_2) = \frac{2}{\|w\|} \quad \leftarrow \text{this is nothing but magnitude of vector.}$$

and

$$y_i = \begin{cases} -1 & w^T x_i + b \leq -1 \\ 1 & w^T x_i + b \geq 1 \end{cases} \quad (\text{label})$$

So max $\left(\frac{2}{\|w\|} \right)$ such that.

$$y_i = \begin{cases} -1 & w^T x_i + b \leq -1 \\ 1 & w^T x_i + b \geq 1 \end{cases}$$

instead of using $\max_{\lambda} \left(\frac{2}{\|w\|} \right)$ we can also try Min which make better score

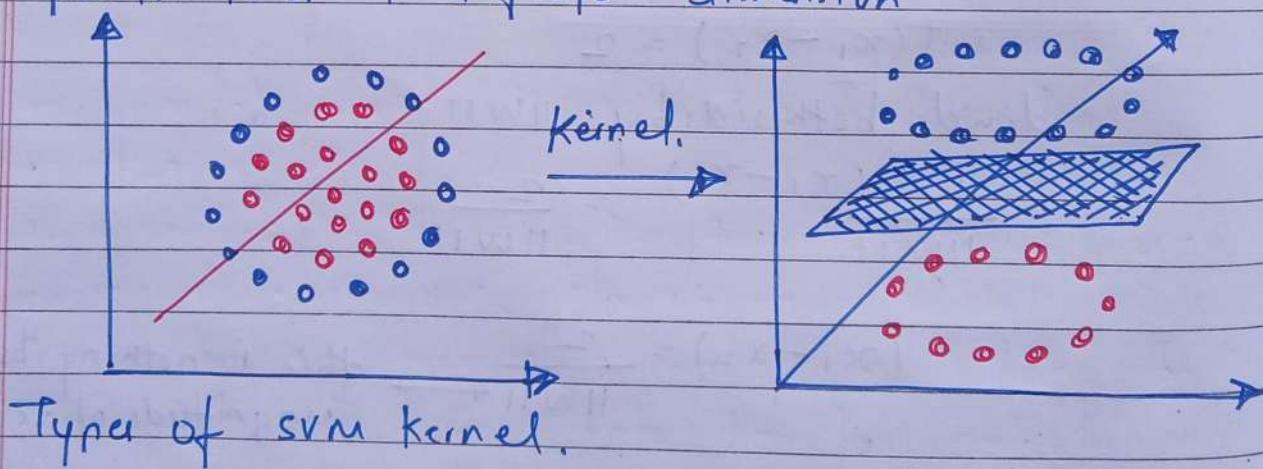
$$\min \left(\frac{\|w\|}{2} \right) + c \times \sum e_i$$

c = Number of error

e_i = Error magnitude

(we all model to train with some error to avoid overfitting i.e. it will be good and train and will bad for test data)

Kernels' in SVM : Generally function of the Kernel is to transform the training set of data so that non-linear decision surface can be transformed to a linear equation in higher number of dimension space it return the inner product between two points in standard feature dimension



- 1) Linear
- 2) polynomial
- 3) Radial Basis function. (rbf)
- 4) sigmoid.

Jupiter's feature (x)

x	-6	-5	-4	-3	-2	-1	0	1	2	3	4	5	6
x^2	36	25	16	9	4	1	0	1	4	9	16	25	36

if you try to plot x on this 1D line.

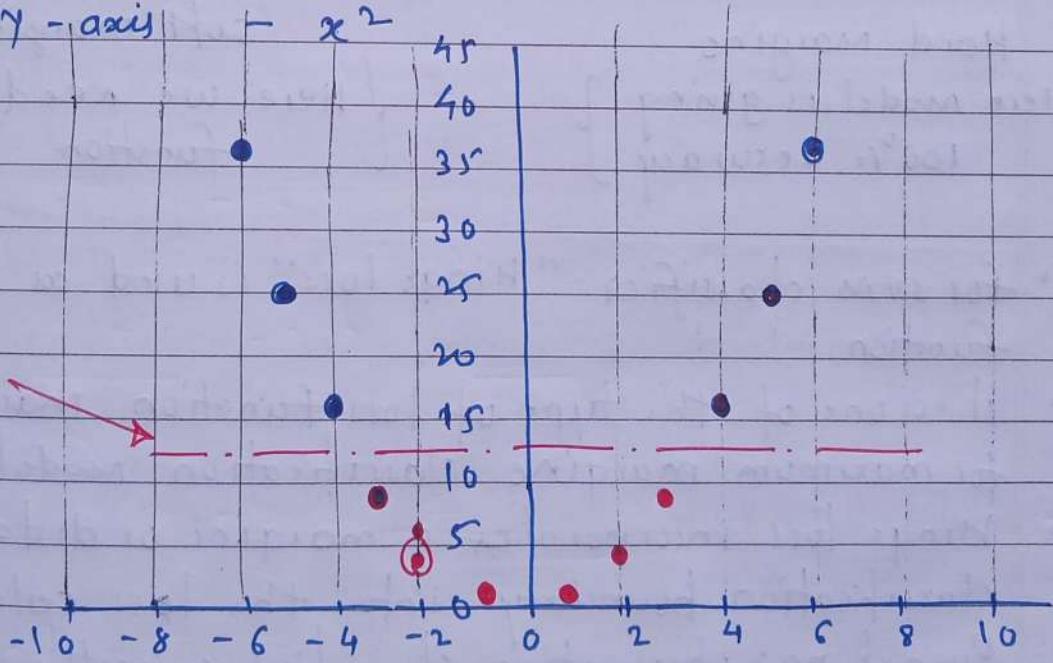


- we can see None of the line could separate the two class perfectly.
- that's why we add another feature which is function of x i.e x^2

x -axis = x

y -axis = x^2

Now this is
separable
data.



1) Linear kernel :- $k(x_1, x_2) = x_1^T \cdot x_2$

{ best suitable for having too many features }

2) polynomial kernel.

$$k(x_1, x_2) = (x_1^T \cdot x_2 + r)^d$$

↑
degree

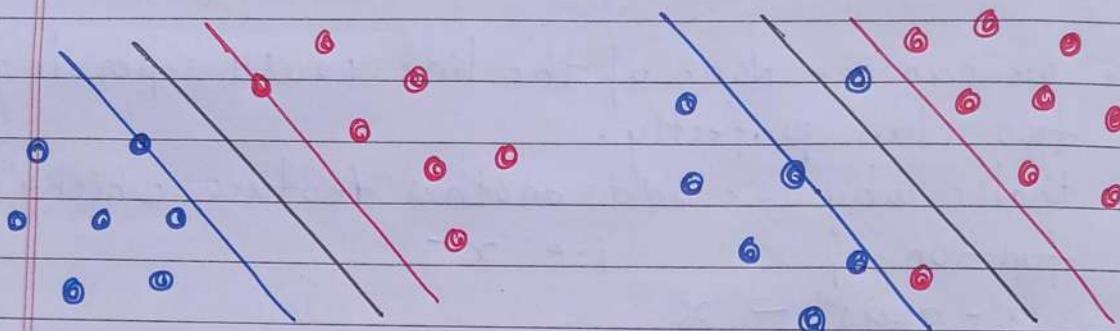
3.) radial basis function. (rbf kernel)

$$k(x_1, x_2) = \exp(-\gamma \cdot \|x_1 - x_2\|^2)$$

4. Sigmoid function

$$k(x_1, x_2) = \tanh(\gamma \cdot x_1 \cdot x_2 + \gamma).$$

Loss function for SVM classifier.



Hard Margin

{ Here model is giving }
100% accuracy }

soft margin

{ Here we need loss }
function }

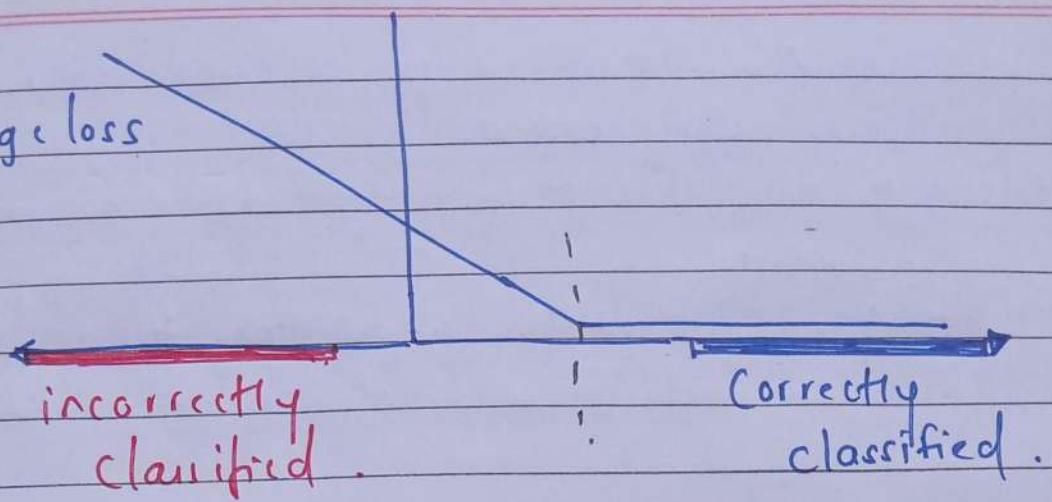
- for SVM classifier "Hinge loss" is used as loss function.
- it is one of the types of loss function mainly used for maximum margin classification model.
- Hinge loss incorporates a margin or distance from classification boundary into the loss calculation. Even if new observation classified correctly they can incur penalty if the margin from decision boundary is not large enough.

$$L = \max(0, 1 - y_i(\omega^T x + b))$$

0 - for correct classification

1 - for wrong classification

hinge loss



wrong
let's talk for miscalcification.

$$y_i = 1, \quad \hat{y}_i = -1 \quad y_i = -1 \quad \hat{y} = 1$$

$$\begin{aligned} L &= (1 - 1)(-1) \\ &= 1 + 1 \\ &= 2 \end{aligned} \quad \begin{aligned} L &= (1 - (-1))(1) \\ &= 1 + 1 \\ &= 2 \end{aligned}$$

{ both are high loss value }

Now talk for correct classification.

$$y_i = 1 \quad \hat{y}_i = 1 \quad y_i = -1 \quad \hat{y}_i = -1$$

$$(0 - (1)(1)) \quad (0 - (-1)(-1))$$

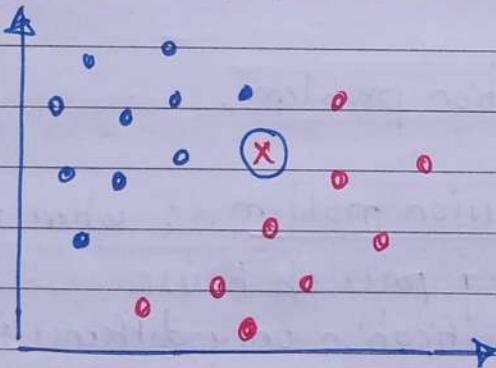
0 - 1	0 - 1
-1	-1

{ both are low loss value }

KNN - K-Nearest Neighbor.

- The abbreviation KNN stands for "K-Nearest Neighbor" it is supervised machine learning algorithm this algorithm can be used to solve both classification as well as regression problem.
- The number of nearest neighbor to a new unknown variable that has to be predicted or classified denoted by the symbol 'K'
- whenever new data will come if new data is close to 1 then prediction will be class 1 otherwise 0 for binary classification same principle for multiclassification problem.
- In general K is odd number

Let understand with some example



Now we have to find new added \textcircled{X}

Working of KNN Algorithm.

Step 1: Loading training as well as test data.

Step 2: Next we choose value of K (hyperparameter) the nearest datapoints. k can be any integer.

Step 3: for each data point (new) test data do the following.

3.1 > calculate distance between test data and each row of training data with the help of distance calculating method like Euclidean, Manhattan or Hamming distance. The most commonly method to calculate is Euclidean.

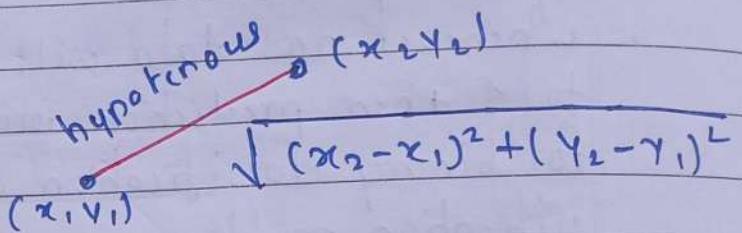
3.2) Now based on its distance value, sort them in ascending order.

3.3) Next it will choose the top k rows from sorted array.

3.4) Now it will assign a class to a test point based on most frequent class of those rows.

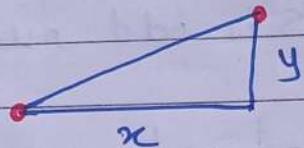
Step 4 : End.

a) Euclidean distance :-



$$\text{hypotenuse} = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

b) Manhattan distance



$$\text{M. distance} = x + y.$$

This is for classification problem.

Now let's see for regression problem. : where output is continuous data. e.g. price of house.

• It is similar to classification only difference is last step where we were considering most frequent class here from k number of output (nearest 'k') we calculate their mean, that's it.

Limitation : Not applicable to huge dataset since calculating distance would consume lot of time.

- Sensitive to outliers
- Inensitive to missing values

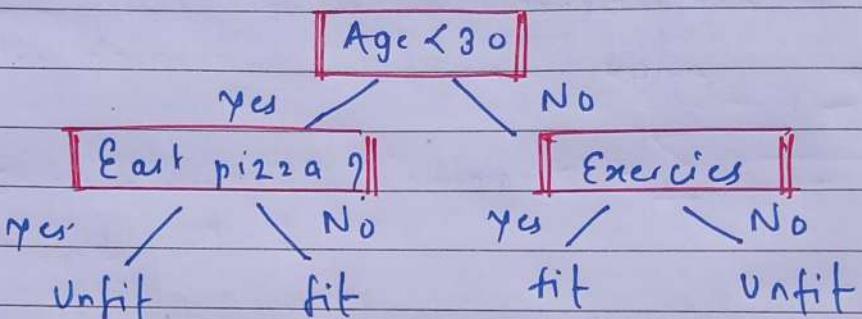
Design Tree.

Show decision Tree.

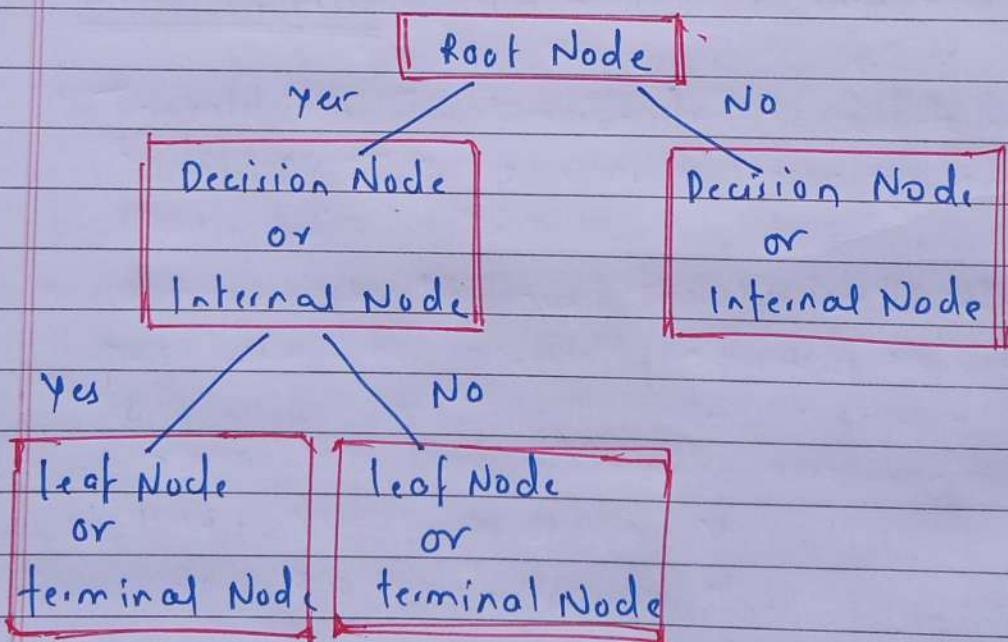
- 1) It is supervised ML Model
- 2) Used both classification & Regression.
- 3) Build Decision Nodes at each step.
- 4) Basis of Tree based model.

Let understand with Example.

is person fit or Not ?



structure & terminology of DT :-



Advantages :-

- 1) Can be used for both classification & Regression
- 2) Easy to interpret
- 3) No need for Normalization or scaling
- 4) Not sensitive to outliers.

